

Variational Methods for Latent Variable Problems (part 2)

Ryan Giordano (for Johns Hopkins Biostats BLAST working group)

Oct, 2021

Massachusetts Institute of Technology

Outline for today:

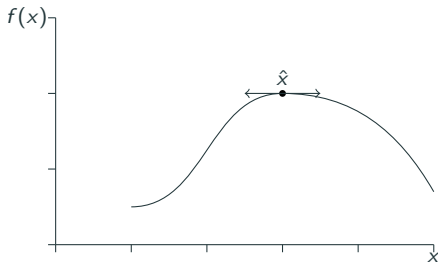
- What counts as variational inference?
- Kullback-Leibler (KL) divergence and “standard” variational inference
- The classical EM algorithm as a special case of variational inference
- Variational inference as a generalization of the EM algorithm
- A quick and incomplete sketch of further topics in variational inference

What counts as variational inference?

Lots of very different procedures go by the name “variational inference.” I propose an (idiosyncratic) encompassing definition based on the use cases and the name:

Variational inference is inference using optimization.

Think “calculus of variations:” an optimum $\hat{x} = \underset{\theta}{\operatorname{argmax}} f(x)$ is characterized by $df/dx|_{\hat{x}} = 0$, i.e. where small variations in \hat{x} result in no changes to the value of $f(\hat{x})$.



By this definition,

- The maximum likelihood estimator (MLE) is VI.
- The Laplace approximation to a Bayesian posterior is VI.
- Markov chain Monte Carlo (MCMC) is not VI.

What counts as variational inference?

A more common definition of VI is the following.

Suppose we have a random variable ξ and a distribution $p(\xi)$ that we want to know.

Let y denote data and θ a parameter. Examples:

- The variable is θ , and we wish to know the posterior $p(\theta|y)$ (Bayes)
- The variable is y , and we wish to know $p(y)$ (MLE)
- The variable is y , and we wish to know the map $\theta \mapsto p(y|\theta) = \int p(y, z|\theta) dz$ (marginal MLE)

Let \mathcal{Q} be some class of distributions which may or may not contain $p(\xi)$.

Variational inference finds the distribution in \mathcal{Q} closest to p according to some measure of “divergence” between distributions:

$$q^* = \operatorname{argmin}_{q \in \mathcal{Q}} D(q, p).$$

The most common choice of “divergence” is the **Kullback-Leibler** (KL) divergence, though other choices are possible (e.g. Li and Turner [2016], Liu and Wang [2016], Ambrogioni et al. [2018]).

KL divergence

The KL divergence is defined as:

$$\text{KL}(q||p) := \mathbb{E}_{q(\xi)} [\log q(\xi)] - \mathbb{E}_{q(\xi)} [\log p(\xi)]$$

Some points to be aware of:

- $\text{KL}(q||p) \geq 0$
- $\text{KL}(q||p) = 0 \Rightarrow p = q$
- $\text{KL}(q||p) \neq \text{KL}(p||q)$
- $\text{KL}(q||p)$ is a “strict” measure of closeness [Gibbs and Su, 2002]

Why use KL divergence?

Phony answer: The KL divergence has an information theoretic interpretation [Kullback and Leibler, 1951].

Real answer: Mathematical convenience (normalizing constants pop out).

Example: the MLE minimizes KL divergence. Suppose that $x_n \stackrel{iid}{\sim} p(\cdot)$, and $q(\cdot|\theta) \in \mathcal{Q}$ is a (possibly misspecified) parameteric family of data distributions. Then

$$\begin{aligned}\hat{\theta} &:= \underset{\theta}{\operatorname{argmin}} \text{KL}(p||q) = \underset{\theta}{\operatorname{argmin}} \left(- \mathbb{E}_{p(x_1)} [\log q(x_1|\theta)] + \mathbb{E}_{p(x_1)} [\log p(x_1)] \right) \\ &= \underset{\theta}{\operatorname{argmax}} \mathbb{E}_{p(x_1)} [\log q(x_1|\theta)] \approx \underset{\theta}{\operatorname{argmax}} \frac{1}{N} \sum_{n=1}^N \log q(x_n|\theta).\end{aligned}$$

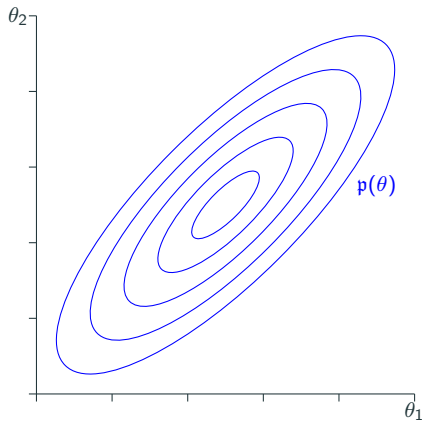
KL divergence exercises

$$\text{KL}(q(\theta) || p(\theta)) = \\ - \mathbb{E}_{q(\theta)} [\log p(\theta)] + \mathbb{E}_{q(\theta)} [\log q(\theta)]$$

$p(\theta)$ = Correlated bivariate normal

$\mathcal{Q} = \{\text{All bivariate normals}\}$

What is $q^*(\theta) = \underset{q \in \mathcal{Q}}{\text{argmin}} \text{KL}(q(\theta) || p(\theta))$?



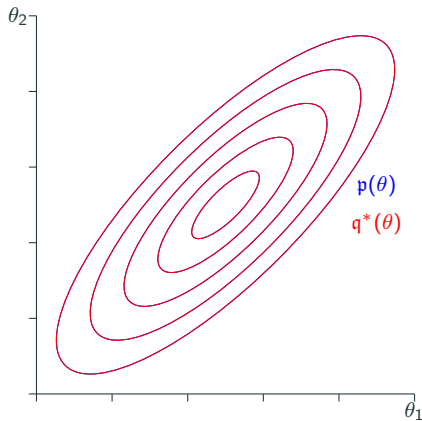
KL divergence exercises

$$\text{KL} (q(\theta)||p(\theta)) = \\ - \mathbb{E}_{q(\theta)} [\log p(\theta)] + \mathbb{E}_{q(\theta)} [\log q(\theta)]$$

$p(\theta)$ = Correlated bivariate normal

$\mathcal{Q} = \{\text{All bivariate normals}\}$

What is $q^*(\theta) = \operatorname{argmin}_{q \in \mathcal{Q}} \text{KL} (q(\theta)||p(\theta))$?



Sufficiently expressive families recover the target distribution.

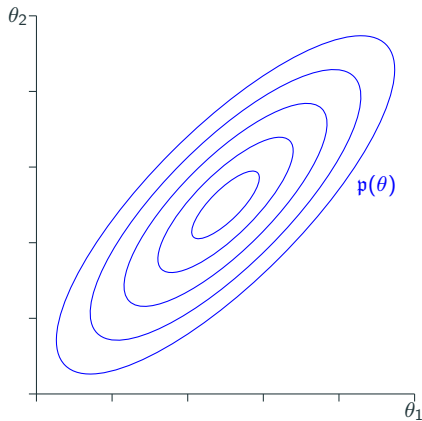
KL divergence exercises

$$\begin{aligned} \text{KL}(q(\theta) || p(\theta)) = \\ - \mathbb{E}_{q(\theta)} [\log p(\theta)] + \mathbb{E}_{q(\theta)} [\log q(\theta)] \end{aligned}$$

$p(\theta)$ = Correlated bivariate normal

$\mathcal{Q} = \{\text{Independent bivariate normals}\}$

What is $q^*(\theta) = \underset{q \in \mathcal{Q}}{\text{argmin}} \text{KL}(q(\theta) || p(\theta))$?



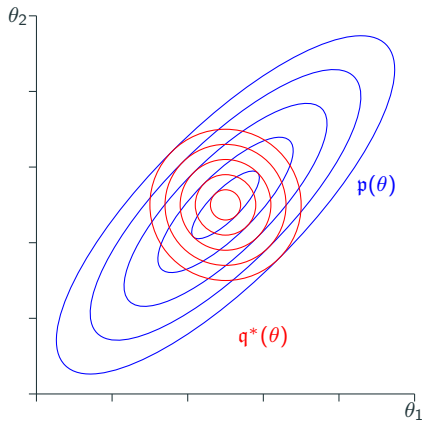
KL divergence exercises

$$\begin{aligned} \text{KL}(q(\theta) || p(\theta)) = \\ - \mathbb{E}_{q(\theta)} [\log p(\theta)] + \mathbb{E}_{q(\theta)} [\log q(\theta)] \end{aligned}$$

$p(\theta)$ = Correlated bivariate normal

$\mathcal{Q} = \{\text{Independent bivariate normals}\}$

What is $q^*(\theta) = \underset{q \in \mathcal{Q}}{\text{argmin}} \text{KL}(q(\theta) || p(\theta))$?



KL minimizers “fit inside” the second argument.

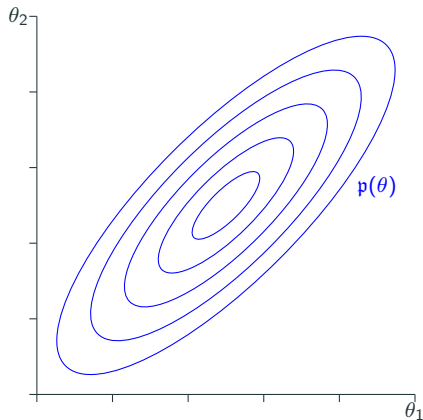
KL divergence exercises

$$\begin{aligned} \text{KL} (p(\theta) || q(\theta)) = \\ - \mathbb{E}_{p(\theta)} [\log q(\theta)] + \mathbb{E}_{p(\theta)} [\log p(\theta)] \end{aligned}$$

$p(\theta)$ = Correlated bivariate normal

$\mathcal{Q} = \{\text{Independent bivariate normals}\}$

What is $q^*(\theta) = \underset{q \in \mathcal{Q}}{\text{argmin}} \text{KL} (p(\theta) || q(\theta))$?



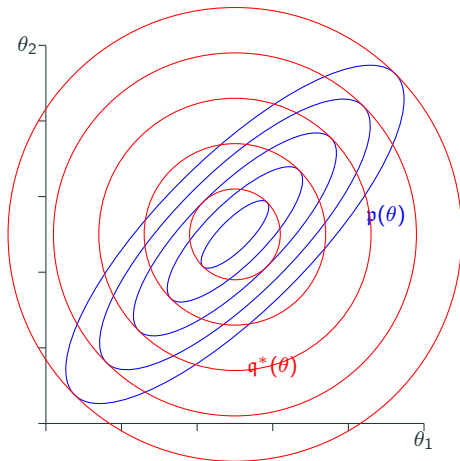
KL divergence exercises

$$\begin{aligned} \text{KL}(p(\theta) || q(\theta)) = \\ - \mathbb{E}_{p(\theta)} [\log q(\theta)] + \mathbb{E}_{p(\theta)} [\log p(\theta)] \end{aligned}$$

$p(\theta)$ = Correlated bivariate normal

$\mathcal{Q} = \{\text{Independent bivariate normals}\}$

What is $q^*(\theta) = \underset{q \in \mathcal{Q}}{\text{argmin}} \text{KL}(p(\theta) || q(\theta))$?



KL minimizers “fit inside” the second argument.

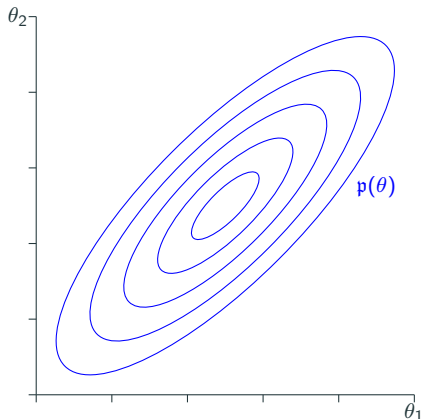
KL divergence exercises

$$\begin{aligned} \text{KL}(q(\theta) || p(\theta)) = \\ - \mathbb{E}_{q(\theta)} [\log p(\theta)] + \mathbb{E}_{q(\theta)} [\log q(\theta)] \end{aligned}$$

$p(\theta)$ = Correlated bivariate normal

$$\mathcal{Q} = \{\text{Bivariate normals}\}$$

What is $q^*(\theta) = \underset{q \in \mathcal{Q}}{\text{argmin}} \left(- \mathbb{E}_{q(\theta)} [\log p(\theta)] + \mathbb{E}_{q(\theta)} [\log q(\theta)] \right)$?



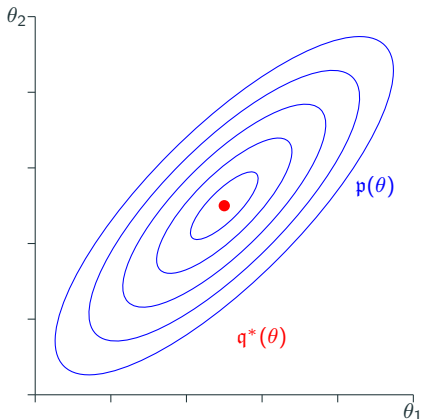
KL divergence exercises

$$\begin{aligned} \text{KL}(q(\theta) || p(\theta)) = \\ - \mathbb{E}_{q(\theta)} [\log p(\theta)] + \mathbb{E}_{q(\theta)} [\log q(\theta)] \end{aligned}$$

$p(\theta)$ = Correlated bivariate normal

$\mathcal{Q} = \{\text{Bivariate normals}\}$

What is $q^*(\theta) = \underset{q \in \mathcal{Q}}{\text{argmin}} \left(- \mathbb{E}_{q(\theta)} [\log p(\theta)] + \mathbb{E}_{q(\theta)} [\log q(\theta)] \right)$?



Without the entropy, the KL minimizer concentrates on the maximum of $\log p(\theta)$.

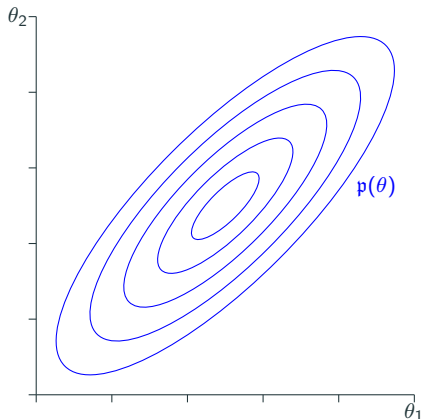
KL divergence exercises

$$\begin{aligned} \text{KL}(q(\theta) || p(\theta)) = \\ - \mathbb{E}_{q(\theta)} [\log p(\theta)] + \mathbb{E}_{q(\theta)} [\log q(\theta)] \end{aligned}$$

$p(\theta)$ = Correlated bivariate normal

$$\mathcal{Q} = \{\text{Bivariate normals}\}$$

What is $q^*(\theta) =$
 $\underset{q \in \mathcal{Q}}{\text{argmin}} \left(- \mathbb{E}_{q(\theta)} [\log p(\theta)] + \mathbb{E}_{q(\theta)} [\log q(\theta)] \right) ?$



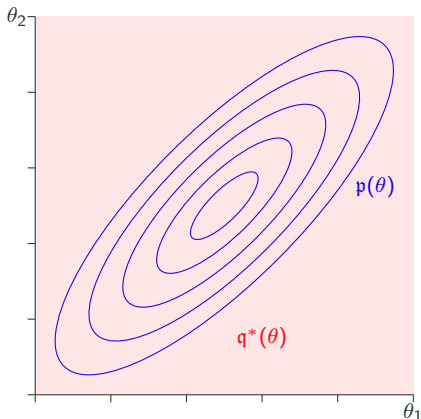
KL divergence exercises

$$\begin{aligned} \text{KL}(q(\theta) || p(\theta)) = \\ - \mathbb{E}_{q(\theta)} [\log p(\theta)] + \mathbb{E}_{q(\theta)} [\log q(\theta)] \end{aligned}$$

$p(\theta)$ = Correlated bivariate normal

$\mathcal{Q} = \{\text{Bivariate normals}\}$

What is $q^*(\theta) =$
 $\underset{q \in \mathcal{Q}}{\text{argmin}} \left(- \mathbb{E}_{q(\theta)} [\log p(\theta)] + \mathbb{E}_{q(\theta)} [\log q(\theta)] \right)$?



Without $\log p(\theta)$, the KL minimizer is infinitely dispersed.

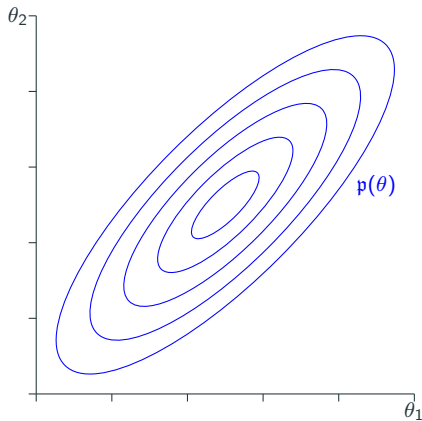
KL divergence exercises

$$\text{KL}(q(\theta) || p(\theta)) = -\mathbb{E}_{q(\theta)} [\log p(\theta)] + \mathbb{E}_{q(\theta)} [\log q(\theta)]$$

$p(\theta)$ = Correlated bivariate normal

$\mathcal{Q} = \{\text{Point masses}\}$

What is $q^*(\theta) = \underset{q \in \mathcal{Q}}{\text{argmin}} \text{KL}(q(\theta) || p(\theta))$?



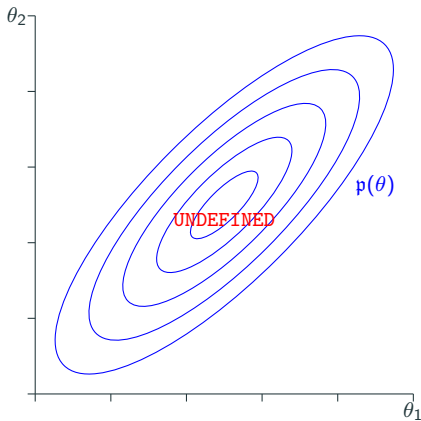
KL divergence exercises

$$\begin{aligned} \text{KL}(q(\theta) || p(\theta)) = \\ - \mathbb{E}_{q(\theta)} [\log p(\theta)] + \mathbb{E}_{q(\theta)} [\log q(\theta)] \end{aligned}$$

$p(\theta)$ = Correlated bivariate normal

$\mathcal{Q} = \{\text{Point masses}\}$

What is $q^*(\theta) = \underset{q \in \mathcal{Q}}{\operatorname{argmin}} \text{KL}(q(\theta) || p(\theta))$?



Without a common dominating measure, the KL divergence is undefined.

KL divergence exercises

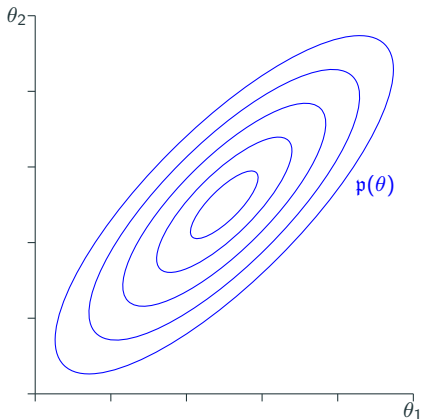
$$\text{KL}(q(\theta) || p(\theta)) =$$

$$- \mathbb{E}_{q(\theta)} [\log p(\theta)] + \mathbb{E}_{q(\theta)} [\log q(\theta)]$$

$p(\theta)$ = Correlated bivariate normal

$\mathcal{Q} = \{\text{BVN with small, fixed variance}\}$

What is $q^*(\theta) = \underset{q \in \mathcal{Q}}{\text{argmin}} \text{KL}(q(\theta) || p(\theta))$?



KL divergence exercises

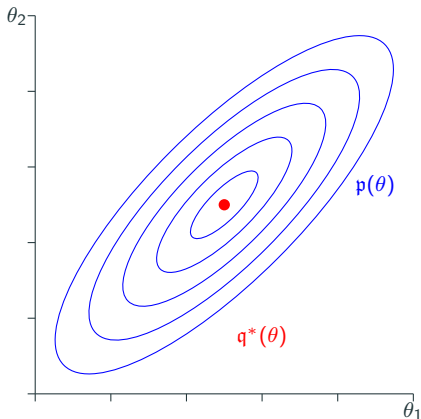
$$\text{KL}(q(\theta) || p(\theta)) =$$

$$- \mathbb{E}_{q(\theta)} [\log p(\theta)] + \mathbb{E}_{q(\theta)} [\log q(\theta)]$$

$p(\theta)$ = Correlated bivariate normal

$\mathcal{Q} = \{\text{BVN with small, fixed variance}\}$

What is $q^*(\theta) = \underset{q \in \mathcal{Q}}{\operatorname{argmin}} \text{KL}(q(\theta) || p(\theta))$?



Sufficiently concentrated distributions with constant entropy act like a point mass at the maximum of $\log p(\theta)$.

- Luca Ambrogioni, Umut Güçlü, Yağmur Güçlütürk, Max Hinne, Eric Maris, and Marcel AJ van Gerven. Wasserstein variational inference. *arXiv preprint arXiv:1805.11284*, 2018.
- Alison L Gibbs and Francis Edward Su. On choosing and bounding probability metrics. *International statistical review*, 70(3):419–435, 2002.
- Solomon Kullback and Richard A Leibler. On information and sufficiency. *The annals of mathematical statistics*, 22(1):79–86, 1951.
- Yingzhen Li and Richard E Turner. Variational inference with rényi divergence. *stat*, 1050:6, 2016.
- Qiang Liu and Dilin Wang. Stein variational gradient descent: A general purpose bayesian inference algorithm. *arXiv preprint arXiv:1608.04471*, 2016.